

AMENDMENTS TO THE CLAIMS

1. (Previously Presented) A server cluster having a plurality of nodes, the server cluster comprising:

a set of base files, wherein said base files are a set of frequently accessed files fitting into a cluster memory of said server cluster, and wherein the base files are logically partitioned into a set of core files having a core size, a set of partitioned files having a partitioned size, and a set of on disk files, and wherein a total of the partitioned size added to the product of the number of said plurality of nodes multiplied by the core size is no greater than the cluster memory; and wherein each node of said plurality of nodes comprises:

at least said set of core files stored locally thereto;

a distributor component for distributing a request to a specific node of said plurality of nodes;

a dispatcher component comprising routing information for said plurality of nodes and replicated across said plurality of nodes, wherein said routing information indicates which said node of said plurality of nodes is for processing said request, said dispatcher component coupled to said distributor component; and

a server component for processing said request, said server component coupled to said dispatcher component;

wherein said plurality of nodes are coupled to a network.

2. (Original) A server cluster as recited in Claim 1 wherein said server cluster is a web server cluster.

3. (Canceled)

4. (Previously Presented) A server cluster as recited in Claim 1 wherein said cluster memory is a combined random access memory of each of said nodes of said server cluster.

5. (Previously Presented) A server cluster as recited in Claim 1 wherein each of said plurality of nodes further comprises a set of partitioned files.

6. (Previously Presented) A server cluster as recited in Claim 1 wherein said set of core files comprises a set of most frequently accessed files of said set of base files.

7. (Previously Presented) A server cluster having a plurality of nodes, wherein each node of said plurality of nodes comprises:

a distributor component for distributing a request to a specific node of said plurality of nodes;

a dispatcher component comprising routing information for said plurality of nodes and replicated across said plurality of nodes, wherein said routing information indicates which said node of said plurality of nodes is for processing said request, said dispatcher component coupled to said distributor component;

a server component for processing said request, said server component coupled to said dispatcher component;

wherein said plurality of nodes are coupled to a network, and wherein each of said plurality of nodes further comprises a set of core files and a set of partitioned files; and

wherein said set of core files is identified by the steps of:

a) logically partitioning said base files into a first subset of files having a first size, a second subset of files having a second size, and a third subset of files having a third size, wherein said base files comprising each of said first subset of files, said second subset of files, and said third subset of files are ordered in decreasing frequency of access;

b) identifying said first subset of files and said second subset of files wherein the total of said second size added to the product of said number of nodes multiplied by said first is less than said cluster memory; and

c) minimizing a total overhead due to the base files wherein said total overhead equals an overhead of said first subset of files plus an overhead of said second subset of files plus said overhead of said third subset of files.

8 – 13. (Cancelled)

14. (Previously Presented) A method for managing request distribution to a set of files stored on a server, said method comprising the steps of:

a) receiving a request for a file at a first node of a plurality of nodes, each of said nodes comprising a distributor component for distributing a request to a specific node of said plurality of nodes, a dispatcher component comprising routing information for said plurality of nodes and replicated across said plurality of nodes, and a server component for processing said request;

b) provided said request is for a core file, processing said request at said first node;

c) provided said request is for a partitioned file, determining whether said request is assigned to be processed by said first node;

c1) provided said request is for a partitioned file assigned to be processed by said first node, processing said request at said first node;

c2) provided said request is for a partitioned file assigned to be processed by another node of said plurality of nodes, forwarding said request to a specific node of said plurality of nodes as indicated by said dispatcher component of said first node and processing said request at said specific node; and

d) provided said request is not for a said core file or a said partitioned file, processing said request at said first node;

wherein each of said plurality of nodes further comprises a set of core files comprising said core file and a set of partitioned files comprising said partitioned file, and wherein said set of core files is identified by the steps of:

a) logically partitioning said base tiles into said set of core files having a core size, said set of partitioned files having a partitioned size, and a set of on disk files having an on disk size, wherein said base files comprising each of said set of core files, said set of partitioned files, and said set of on disk files are ordered in decreasing frequency of access;

b) identifying said set of core files and said set of partitioned files wherein the total of said partitioned size added to the product of said number of nodes multiplied by said core size is less than said cluster memory and

c) minimizing a total overhead due to the base files wherein said total overhead equals an overhead of said core set of files plus an overhead of said partitioned set of files plus said overhead of said on disk set of files.

15. (Previously Presented) A method for identifying a set of frequently accessed files on a server cluster comprising a number of nodes, said method comprising the steps of:

- a) defining a set of base files, wherein said base files are a set of frequently accessed files fitting into the cluster memory of said server cluster, said base files ordered in decreasing frequency of access;
- b) logically partitioning said base files into a first subset of files having a first size, a second subset of files having a second size, and a third subset of files having a third size, wherein said base files comprising each of said first subset of files, said second subset of files, and said third subset of files are ordered in decreasing frequency of access;
- c) identifying said first subset of files and said second subset of files wherein the total of said second size added to the product of said number of nodes multiplied by said first size is not greater than said cluster memory; and
- d) minimizing a total overhead due to the base files wherein said total overhead equals an overhead of said first subset of files plus an overhead of said second subset of files plus said overhead of said third subset of files.

16. (Original) A method as recited in Claim 15 wherein said server cluster is a web server cluster.

17. (Original) A method as recited in Claim 15 wherein said cluster memory is a combined random access memory of each of said nodes of said web server cluster.

18. (Original) A method as recited in Claim 15 wherein said first subset of files is a set of core files and wherein said first size is a core size.

19. (Original) A method as recited in Claim 15 wherein said second subset of files is a set of partitioned files and wherein said second size is a partitioned size.

20. (Original) A method for determining a set of $Files_{core}$, said method comprising the steps of:

- a) defining a set of $BaseFiles$ as a set of frequently accessed files fitting into a $ClusterRAM$, said $BaseFiles$ ordered in decreasing frequency of access;
- b) logically partitioning said $BaseFiles$ into a $Files_{part}$, a $Files_{core}$ and a $Files_{on\ disk}$ wherein $BaseFiles = Files_{part} + Files_{core} + Files_{on\ disk}$;
- c) identifying said set $Files_{part}$ and said set $Files_{core}$ according to $N \propto Size_{core} + Size_{part} \leq ClusterRAM$; and
- d) minimizing $OH_{BaseFiles}$ according to $OH_{BaseFiles} = OH_{part} + OH_{core} + OH_{on\ disk}$.

21. (Previously Presented) A server cluster as recited in Claim 1 wherein each file of said set of partitioned files is assigned to a particular one of said plurality of nodes for processing thereof.

22. (Previously Presented) A server cluster as recited in Claim 21 wherein said routing information indicates which of said plurality of nodes is assigned for processing each of said set of partitioned files.

23. (Previously Presented) A server cluster as recited in Claim 1 wherein said dispatcher component of a node that receives said request (“the receiving node”) comprises logic for determining whether said request is a request to access a core file, and if determined that said request is a request to access said core file then said dispatcher component of the receiving node assigns the server component of the receiving node to process said request.

24. (Previously Presented) A server cluster as recited in Claim 23 wherein said dispatcher component of the receiving node further comprises logic for determining whether said request is a request to access a partitioned file, and if determined that said request is a request to access said partitioned file then said dispatcher component of the receiving node determines one of said plurality of node assigned to process the requested partitioned file.

25. (Previously Presented) A server cluster as recited in Claim 24 wherein if said dispatcher component determines that the receiving node is assigned to process the requested partitioned file, the dispatcher component assigns the server component of the receiving node to process said request.